



HAL
open science

Simplified pangenome graph traversals with PSSM scoring: search for motifs differentials

Julien Guidihounme, Simon De Givry, Benjamin Linard

► To cite this version:

Julien Guidihounme, Simon De Givry, Benjamin Linard. Simplified pangenome graph traversals with PSSM scoring: search for motifs differentials. JOBIM 2024, Jun 2024, Toulouse, France. hal-04656256

HAL Id: hal-04656256

<https://hal.science/hal-04656256v1>

Submitted on 22 Jul 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0 International License

Simplified pangenome graph traversals with PSSM scoring : search for motifs differentials

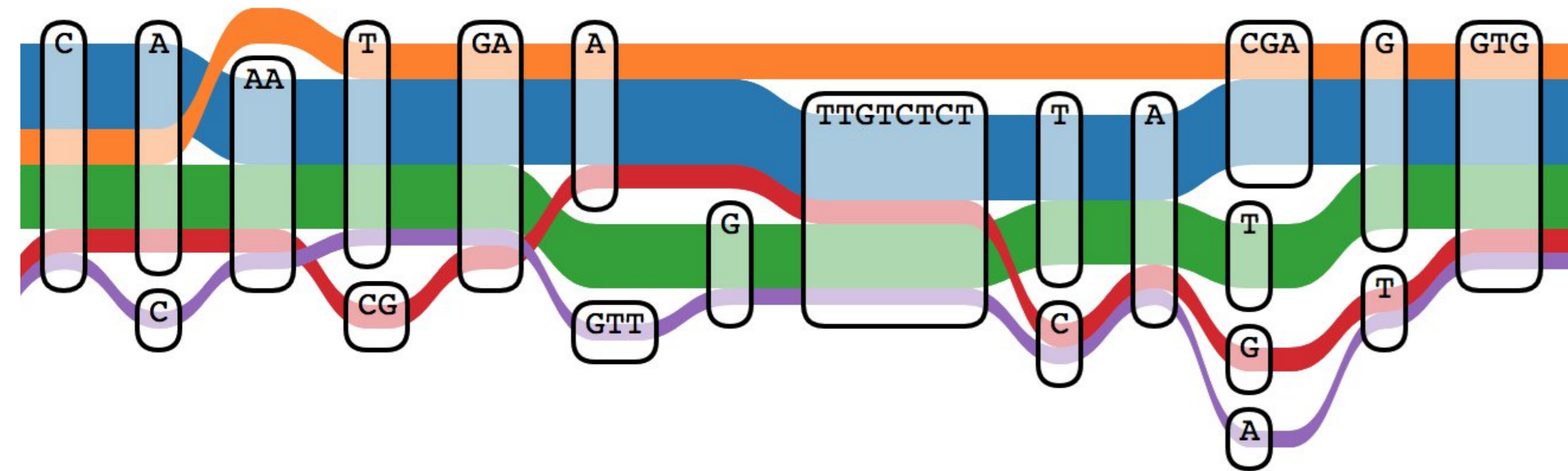
Julien GUIDIHOUNME¹, Simon DE GIVRY¹ and Benjamin LINARD¹

¹ Unité de Mathématiques et Informatique Appliquées, INRAE, Castanet-Tolosan, France
Corresponding Author: julien.guidihounme@inrae.fr

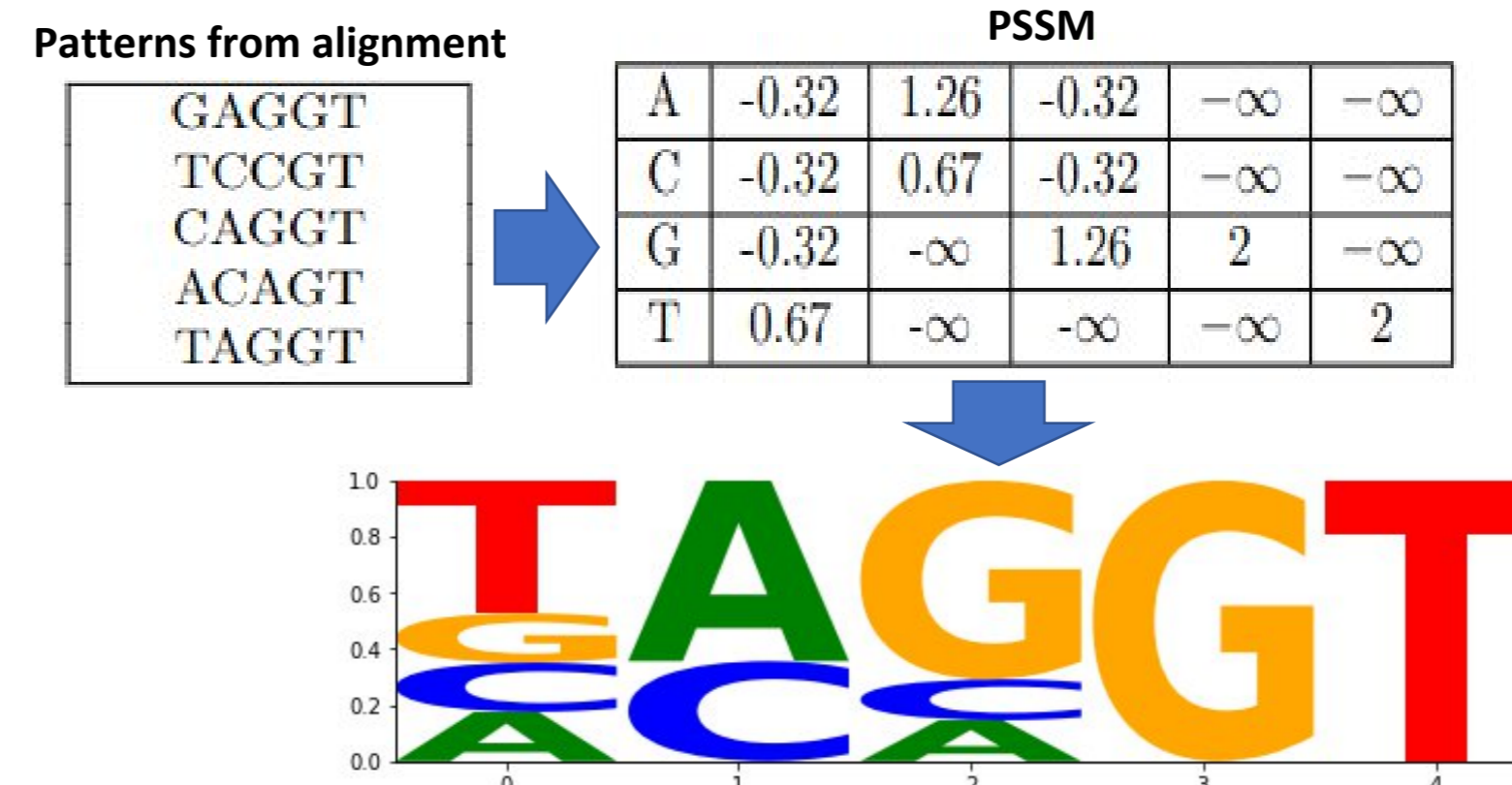
CONTEXT

A pangenome Variation Graph (VG) is a data structure that integrates complete information about genetic variations in a species or group of species [1]. Computational pan-genomic is beginning to be discussed as a viable alternative to many standard approaches previously based on sequences [2]. Although in itself a compression of genome sequences, a fundamental problem is the size of this graph (tens to hundreds of millions of nodes): the application of numerous algorithms based on traversals is not possible and heuristics and indexing techniques are necessary to be able to manipulate and query it. We are interested in developing an approach to search for functional differentials in a VG via the detection of sequence pattern changes between the genomes included in the graph. To do so, we search for probabilistic pattern described via a PSSM (Position Specific Scoring Matrix)[3]. Our approach aims to output all significant PSSM score differentials that can be extracted from a VG. Note we detect motifs on node boundaries only (the score difference is either null or undefined for motifs inside nodes).

Overview of a pangenome graph [4]



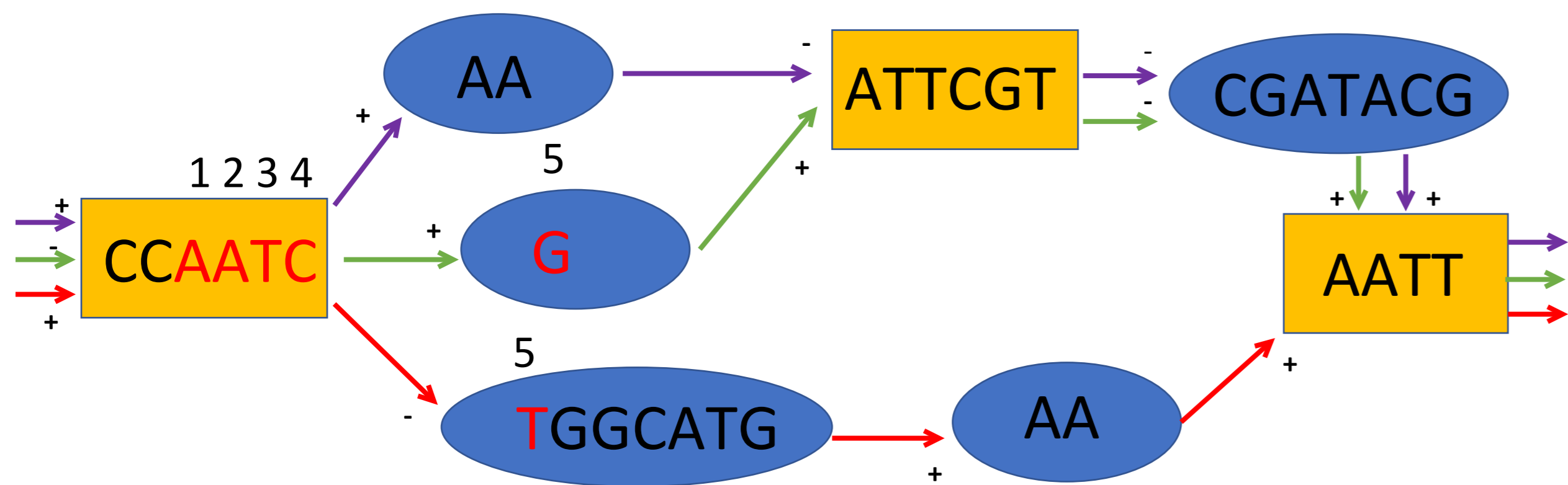
Functional motifs described via a PSSM :



- Determine the score associated with the motif (PSSM) in the graph
- Find for the differential within score the genomic species integrated into the graph
- Scale problem: whence our method on subgraphs defined by pair of genomes to find the differential score

METHOD OVERVIEW

- Our approach iterates on genome **pairs**
- Below example shows 3 genomes (3 possible pairs)



- Yellow nodes** : for each pair, the "bifurcation" nodes are listed, they are the starting point to define sequence positions allowing a putative differential
- In red text**: for the pair (green, red) and for a PSSM matrix of length M, positions where scoring is necessary are selected.
- Scoring**: We use classic PSSM scoring algorithms[5][6] to detect a differential. Positions in the bifurcation node are computed once.
- Iteration**: this process is repeated for all bifurcations, and all genome pairs

FORMALISM & IMPLEMENTATION

Definition : Intersection of genomes $I : \mathcal{P}(\mathcal{P}) \rightarrow \mathcal{P}(\mathcal{V}), J \mapsto \bigcup_{v \in \bigcap_{p \in J} p} \{v\}$

Property :

Let $x \in I(\{p_i, p_j\}), i \neq j, (i, j) \in \{1, 2, \dots, |\mathcal{P}|\}^2$
 $p_i = (v_0, v_1, \dots, v_{n-1}, x, v_{n+1}, \dots, v_{|p_i|})$ and $p_j = (s_0, s_1, \dots, s_{m-1}, x, s_{m+1}, \dots, s_{|p_j|})$
 such that $f_i = (v_0^i, v_1^i, \dots, v_{n-1}^i, x^i, v_{n+1}^i, \dots, v_{|f_i|}^i)$ and $f_j = (s_0^j, s_1^j, \dots, s_{m-1}^j, x^j, s_{m+1}^j, \dots, s_{|f_j|}^j)$
 with $i_j \in \{+, -\}, l \in \{0, 1, \dots, |f_l|\}$ and $j_t \in \{+, -\}, t \in \{0, 1, \dots, |f_t|\}$.

Then x is a **bifurcation node** if it satisfies one of the two conditions :

- v_{n-1} et s_{m-1} exist, then $v_{n-1} \neq s_{m-1}$ ou $i_{n-1} \neq j_{m-1}$
- v_{n+1} et s_{m+1} exist, then $v_{n+1} \neq s_{m+1}$ ou $i_{n+1} \neq j_{m+1}$

Implementation: C++ code, command-line tool

Input:

pangenome graph in GFA format + some PSSM matrix in jasper format

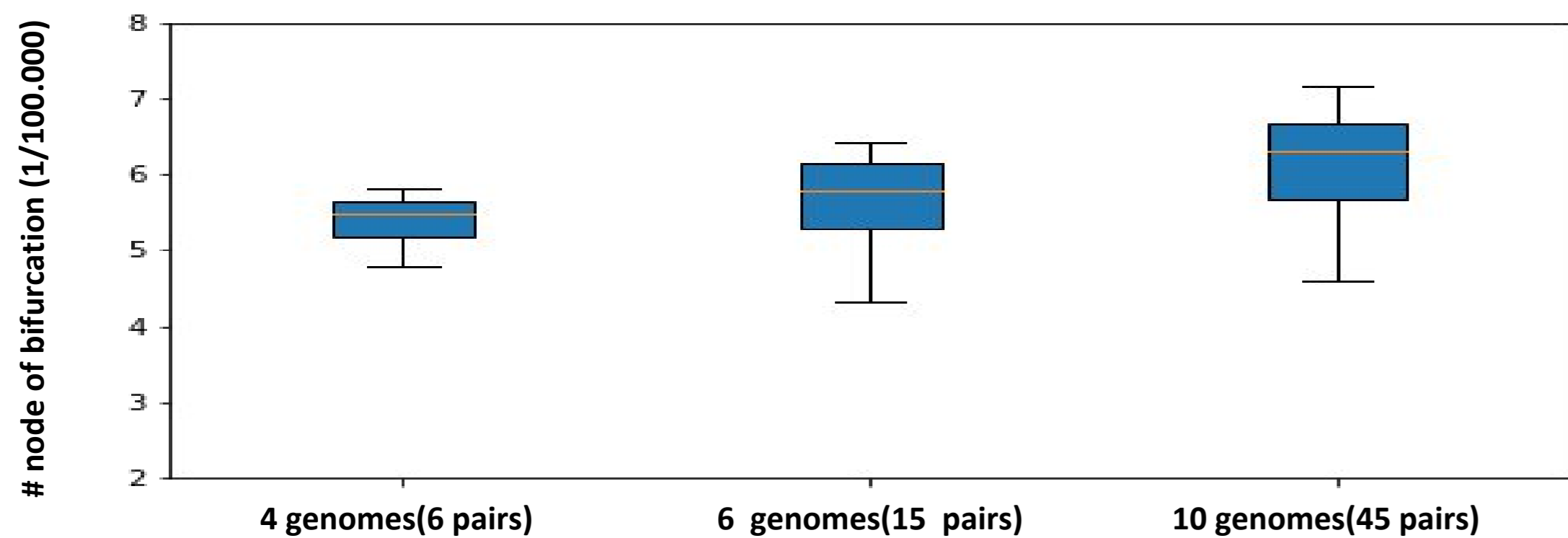
Output:

a list which contains differentials: the score differential value, the index of the bifurcation node in the graph, the corresponding genome positions

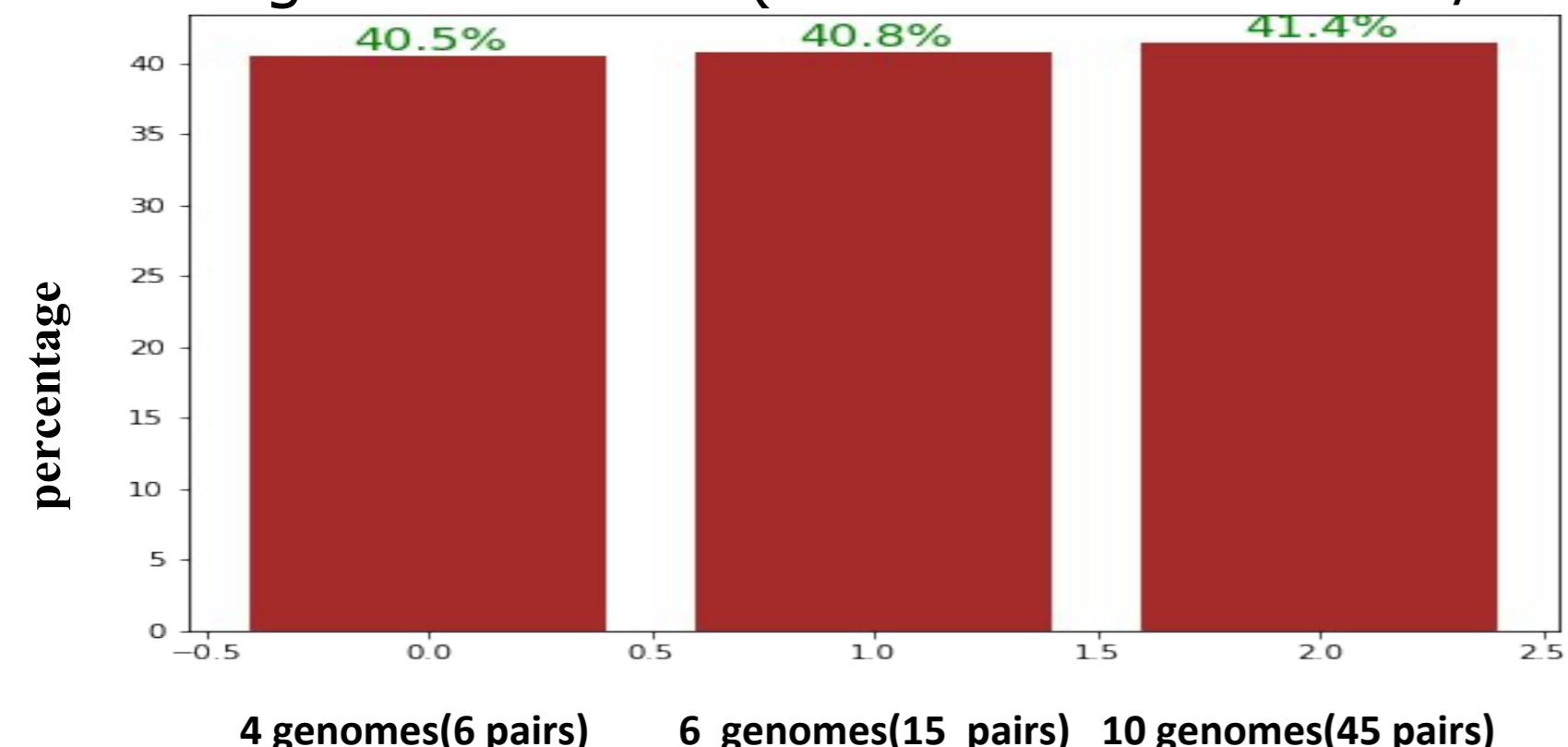
PRELIMINARY RESULTS

Tendency of bifurcation node on the VG of cabbage with 4, 6 and 10 genomes for the same chromosome. (Data courtesy of Fabrice Legai & Alexis Merguez, INRAE, Génotoul Bioinfo)

node of bifurcation stable with the growth of the graph

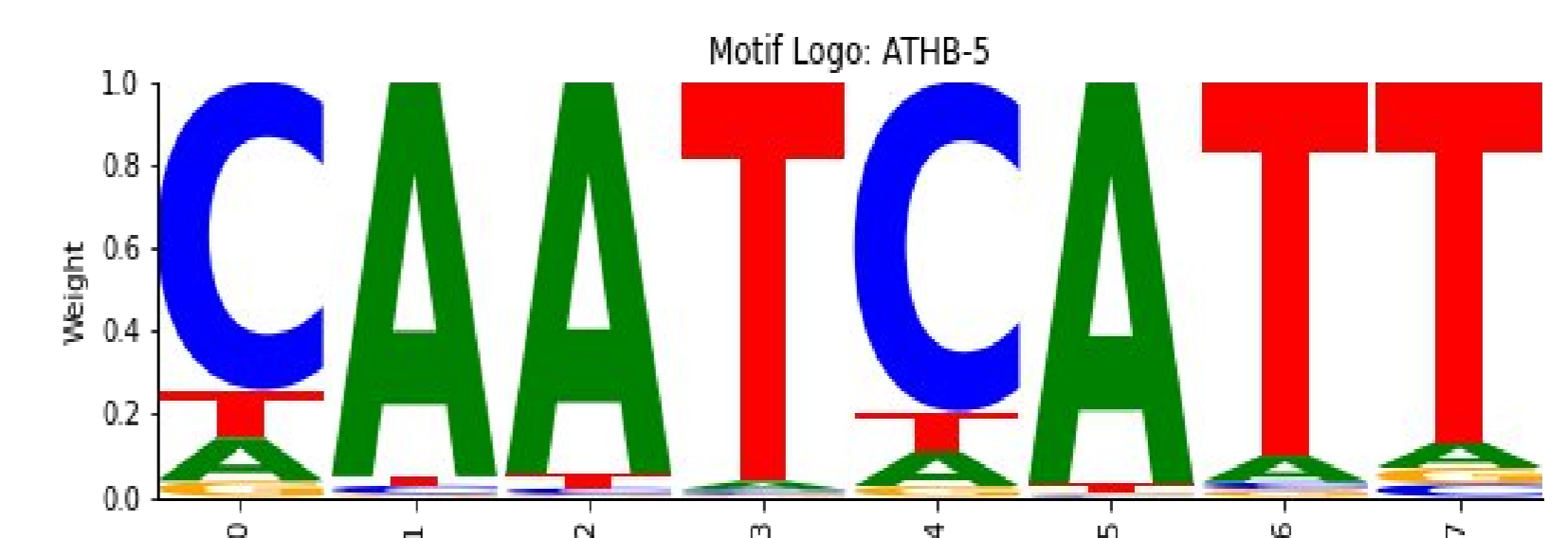


Percentage of the ratio (# node of bifurcation/# node total)



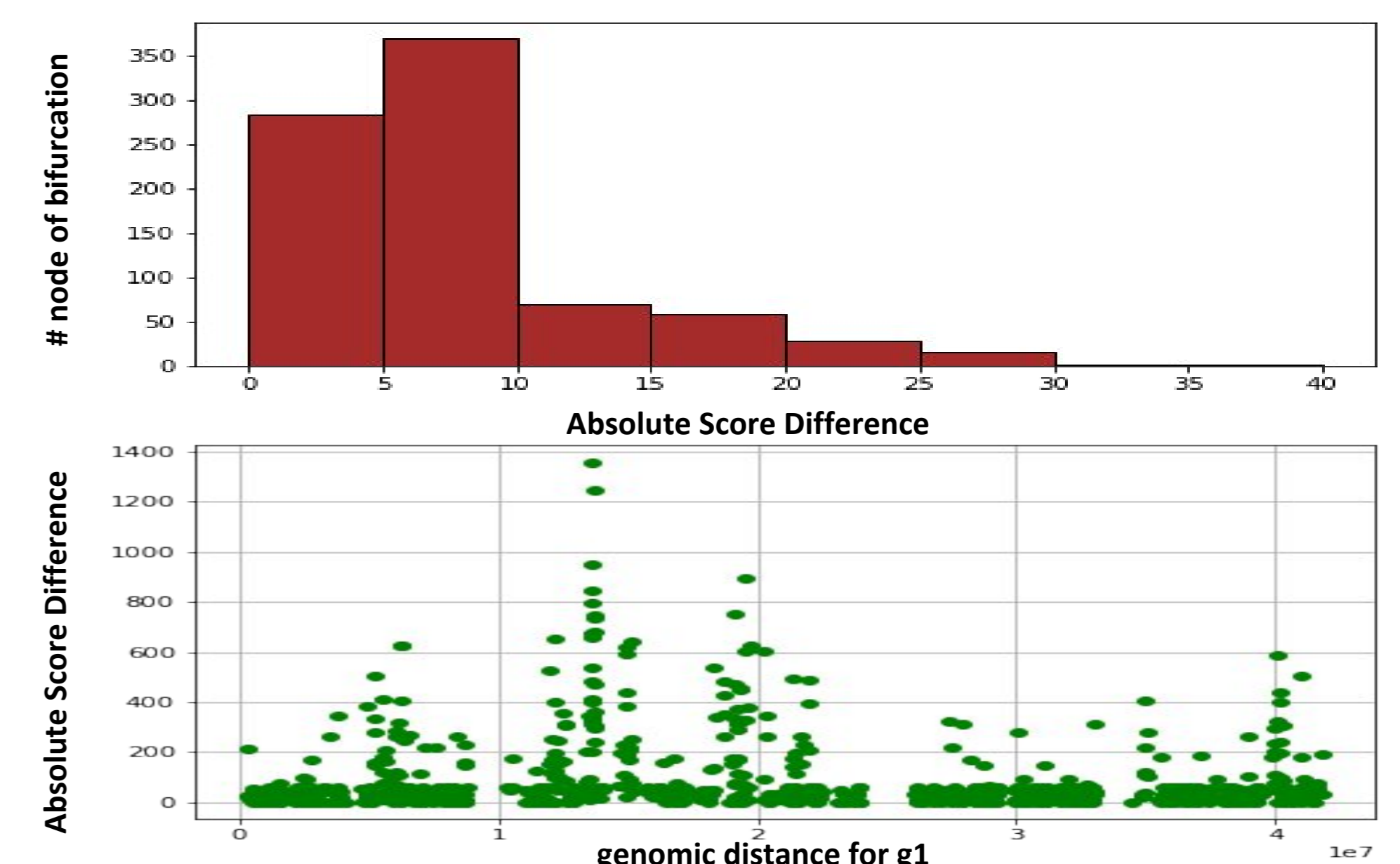
Score differentials for the ATHB-5 motif in an apricot graph : 6 genomes (genome size ~230mb), chromosome 1 only, 2210860 nodes and 3026030 edges. (Data Courtesy of Véronique Decroocq's team, INRAE, BFP, Bordeaux)

ATHB5: a regulatory motif from the jasper database, with a role in the growth and germination of seeds (BL Raminger et al, 2023).



- Score differential for two genomes g_1 and g_2 in the apricot graph by setting the p-value of g_1 to 10^{-4} .

- 350 regions (bifurcation nodes) show a low score differential (interval [0-10]).
- Less than 50 regions show a differential >10, e.g. a significant motif change with putative change/loss of function.



DISCUSSION

- For a VG of one apricot chromosome, we detected 1029524 regions where we calculate score differentials (all 15 pairs of genome) in then less than 15 minutes.
- We continue to optimize the algorithm to avoid redundant computations : use of a score threshold, reuse of score computation between pairs ...
- Code needs to be extended to allow as inputs a collection of PSSM matrices.